

# Are Recent Terrorism Trends Reflected in Social Media?

Ivana Terziyska  
College of Engineering  
Cornell University  
Ithaca, NY, USA  
Email: it76@cornell.edu

Setu Shah  
School of Engineering and Technology  
IUPUI  
Indianapolis, IN, USA  
Email: setshah@iupui.edu

Xiao Luo  
School of Engineering and Technology  
IUPUI  
Indianapolis, IN, USA  
Email: luo25@iupui.edu

**Abstract**—Social media plays an important role in shaping the beliefs and sentiments of an audience regarding an event. A comparison between public data sets that have holistic features and social media data set that include more user features would give insight into the spread of misinformation and aspects of events that are reflected in user behavior. In this research, we compare the trends identified in the public data set - Global Terrorism Database (GTD) with the trends reflected through the social media data obtained using the Twitter API. The unsupervised learning algorithm Self-Organizing Map (SOM) is used to identify the features and trends summarized by the clusters. The results show discrepancies in the features and related trends of terrorism events in the GTD data set and obtained Twitter data set to suggest some media bias and public perception on terrorism.

**Index Terms**—Unsupervised learning, data visualization, social media, terrorism

## I. INTRODUCTION

Terrorism is an increasingly expansive threat to global security and is a continuous and complex focus of research. There is significant discourse about the way terrorism uses media and how terrorism is represented in media. For instance, certain events are deemed “underreported by the media [1], and terrorist organizations have been found to use social media as a digital strategy to promote terror [2]. Consequently, social media platforms such as Facebook are designing weighting algorithms to detect and censor terrorist activity, propaganda, and intentional misinformation [3]. Although comprehensive data exists until 2015 on terrorist attacks, it does not portray the constantly changing and potentially biased social media user interest in current topics. Public awareness, understanding, and recognition of current events is critical towards developing and supporting policies to mitigate their effects.

Data mining using the social media platform Twitter has emerged as a popular tool to extract large volumes of data to gain insight into user behavior, with applications ranging from event classification, disaster and disease surveillance, and opinion mining and event prediction [6] [7] [8]. Our research proposes a novel method of using machine learning to analyze the GTD data and Twitter, and to visualize and infer trends.

There are three major components of our research. First, in-depth analysis on recent terrorist trends is performed using an SOM to determine clusters of similar events in recent history.

Then, the weight of each feature is compared to identify the most related and influential features of terrorist attacks. Finally, SOM is performed to generalize the trends of terrorist activity by identifying the most related features over four years. We propose a method of comparison between current focus on terrorism in social media and pre-existing data by collecting tweets using the Twitter REST API. We classify this information in accordance with the features discovered in the first component. The SOM is employed to partition the data into clusters and compare it to the existing database, as well as identify current trends in social media data related to terrorism.

The remainder of the paper is organized as follows: Section II introduces related work using similar methods and databases and explains the novelty of this particular research. Section III describes the proposed methodology and the self-organizing map. The data sets, specific data pre-processing, and feature extraction, experimental results and discussions on the findings are detailed in Section IV. Section V concludes the paper and proposes future work to elaborate and improve upon this research.

## II. RELATED WORK

Previously, researchers have used different supervised and unsupervised learning algorithms on the GTD data set to identify outliers, classify events and predict terrorism event. Meng et al. [4] implemented both supervised and unsupervised machine learning methods using the GTD data set. They used a modified K-Means clustering algorithm and defined parameters such as cluster size and radius to determine outlier instances in the GTD data set. Their objectives were to decrease human error rate in creating the database and to determine significant events. Supervised classification methods such as Naive Bayes, Support Vector Mechanism, and Logistic Regression were used to test the effectiveness of this classification by using subsets of the GTD for training and testing data. Research by Misra et al. [5] examined neural networks and logistic regression as mechanisms to identify and predict terrorist attacks based on input features and to reduce classification error rate.

Data mining algorithms applied to Twitter data have been employed for terrorism event classification, analysis, and pre-

diction. Oh et. al [9] presented a model based on Situation Awareness (SA) theory to analyze the content of Twitter postings of the 2008 Mumbai terror attacks, examine the use of Twitter as a participatory emergency reporting system, and propose a conceptual framework for analyzing information control and reports. Cheong and Lee [10] chronicled user sentiment and response to terrorist attacks using Twitter data to create graphical visualizations of information, extract knowledge from Twitter messages, and reveal potential responses to terrorist threats.

To the best of authors' knowledge, the proposed research is the first to compare and visualize trends identified through GTD data set against the trends identified by analyzing the Twitter data set through employing the unsupervised machine learning algorithm SOM.

### III. PROPOSED METHODOLOGY

In order to compare the trends identified through GTD data set and Twitter data, pre-processing, and feature extraction/selection have been used before the unsupervised learning algorithm SOM is used to cluster and visualize the trends. Figure 1 demonstrates an overview of the proposed data analytics process. Indeed, different data pre-processing and feature selection steps have been used on GTD and Twitter data respectively. The following sub-sections provide the details of these steps and a brief description of the SOM algorithm.

#### A. Pre-processing

There are total of 45 features in the GTD data set; however, 28 numerically represented features were investigated in this research. Data instances that have missing values of any of the 28 features were excluded from the analysis. The feature representing country name was replaced with the country's respective Global Terrorist Index or GTI, to give a more comprehensive label to each country.

The Twitter REST API [14] was used to extract tweets by using the query keywords "terrorism" and "terrorist". Tweets that included the keywords but did not mention terrorism as defined by the GTD, as well as repeated tweets not classified as retweets, were not considered.

#### B. Feature Extraction

In order to further identify the features that mostly contribute the patterns or clusters within the GTD data set, the selected 28 numeric features were used to train an SOM. Based on the analysis of the input weights of each feature and the data that hit to the neurons with most hits on the SOM, 9 features that include GTI (the Global Terrorist Index), region ID (the geographic ID), crit1 (the motive of the attack), crit3 (whether the attack was outside humanitarian law), success (whether the attack was successful), the target type, the target subtype, the weapon type, and property (if there was property damage or not) were identified as the features that contributed the most to the cluster distributions. From further analysis on these features based on all data instances, 4 features that are

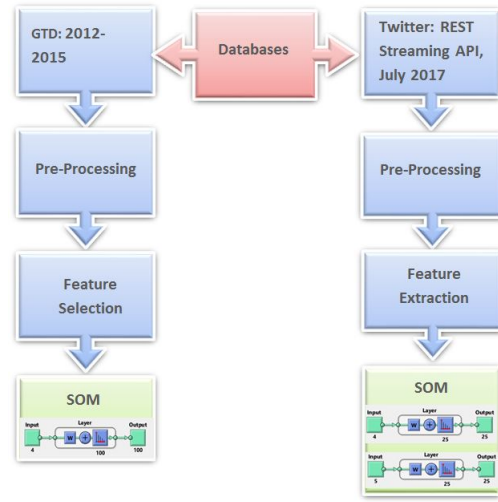


Fig. 1. Overview of the Proposed Data Analytics Process

common to entire data set and can also be extracted from the Twitter data were selected. The details of the features are given in Table I.

In order to extract the corresponding features to the GTD data set from the tweets, the location, tweet text, followers, retweets, favorites, and language were extracted and stored. Only tweets that contained a location can be matched to a real country were collected. For this research, only tweets in English were considered for text-based features, comprising 94% of the total corpus. A manual reviewing process was employed to go through all the tweets to obtain features from the tweets that corresponded to the features of GTD data set. The feature crit1 (motive for attack) is assigned as value 1 if a specific political, economic, social, or religious motive is mentioned in the tweet as a motive of the attack or terrorism in general. Target type is also manually assigned. If a specific target is mentioned in the tweet, the corresponding code in the GTD codebook is assigned. Otherwise, 0 is assigned. Other than the 4 features that match the GTD data set that are extracted from tweets, a sentiment score and favorite score for each tweet are also extracted. A sentiment score (-5 to 5)[15] was assigned to each tweet based on a pre-existing open source library.

#### C. Self-Organizing Map and Data Visualization

The algorithm responsible for the formation of the SOM [11] involves three basic steps after initialization: sampling, similarity matching, and updating. These three steps are repeated until the formation of the feature map has completed. Each neuron  $i$  has a  $d$ -dimensional prototype weight vector  $W_i = W_{i1}, W_{i2}, \dots, W_{id}$ . Given  $X$  is a  $d$ -dimensional sample data (input vector), the algorithm is summarized as follows:

- Initialization:

Choose random values to initialize all the neuron weight vectors  $W_i(0), i = 1, 2, \dots, M$ , where  $M$  is the total number of neurons in the map.

TABLE I  
REPRESENTATION OF FEATURES IN DATASETS

Feature Name	GTD Feature Description	Twitter Feature Description	Measure
GTI	Measure of Terrorist Activity	GTI of User Location	0-10 scale
Region ID	Geographic Region	Region ID of User Location	1-12 Category
crit1	motive for attack(religious, political, social, economics )	motive mentioned in Tweet text	Binary
targtype	target type (e.g government, military, etc)	target type mentioned in Tweet	1-22 Category
favorites	_____	popularity of tweet	number of favorites
sentiment score	_____	user sentiment	-5 to 5

- Sampling:  
Draw a sample data  $X$  from the input space with a uniform probability.
- Similarity Matching:  
Find the best matching unit (BMU) or winner neuron of  $X$ , denoted here by  $b$  which is the closest neuron (map unit) to  $X$  in the criterion of minimum Euclidean distance, at time step  $n$  ( $n^{th}$  training iteration).

$$b = \arg \min_i \|X - W_i(n)\|, i = 1, 2, \dots, M \quad (1)$$

- Updating:  
Adjust the weight vectors of all neurons by using an update formula, so that the best matching unit (BMU) and its topological neighbors are moved closer to the input vector  $X$  in the input space.
- Continuation:  
Continue with sampling until no noticeable changes in the feature map are observed or the pre-defined maximum number of iterations is reached.

The most commonly used visualization techniques of SOM are the U-Matrix and Hit histogram. The U-matrix holds all distances between neurons and their immediate neighbor neurons. It gives a direct visualization of the number of clusters and their distribution on a two dimensional space. Each input data instance in the data set can be projected (hit) to the closest neuron on a trained SOM map. The hit histogram is constructed by counting the number of hits each neuron receives from the input data set. On the hit histogram, the larger the shaded area is on the neuron, the more hits the neuron receives.

#### IV. EXPERIMENTS AND RESULTS

##### A. Data sets

The Global Terrorism Database (GTD)[12] is collected by researchers at the University of Maryland and made available through the National Center for the Study of Terrorism and Responses to Terrorism (START). This database includes over 170,000 data points representing attack records with 45 features from 1970-2015. For this experiment, the most recent GTD data was used for the years 2012-2015, the most recent years available. The Global Terrorist Index (GTI) database was used in place of the name of each country in the GTD [13]. By using the GTI database, each country is assigned a value that represents “lives lost, injuries, property damage

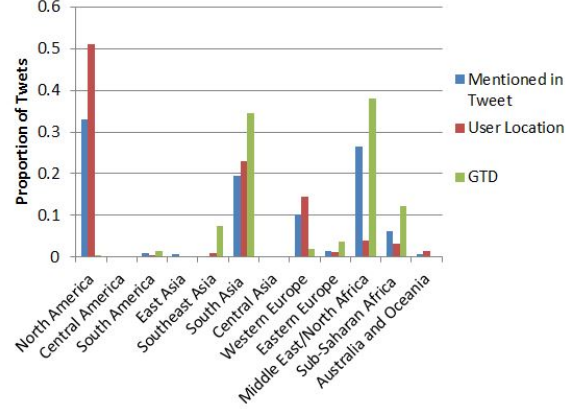


Fig. 2. Regional Representation of the GTD Dataset

and the physiological after effects of terrorism”. After data pre-processing and feature selection, this research uses 51988 data instances of the GTD data set.

Through the use of the Twitter API, the most recent tweets between July 11 and July 18, 2017 were extracted. After pre-processing, there are 1785 tweets used for this research. Among these tweets, 28% mentioned motive for attack, and 13% mentioned a particular attack type. Figure 2 shows the distribution of particular regions that were mentioned in the tweets.

##### B. Clustering Results, Interpretation and Discussion

Self-Organizing maps were trained by using the neural network function in Matlab on each data set. A 10 by 10 map was trained for the GTD, and a 5 by 5 map was trained for Twitter data with and without the sentiment and favorites feature.

In this research, 10000 iterations were used to ensure that the SOM converged to present the cluster distributions of the data. The U-Matrix and hit histogram were used to visually demonstrate the clustering and relationships between features of the data. Figure 3 to 6 present the U-Matrices and hit histograms of four trained SOMs based on four different data sets: GTD data set with four features, Twitter data set with four features, Twitter data set with ‘favorite’ as additional feature and Twitter data set with ‘sentiment’ as additional feature. The left side of the figures are U-matrices which represent distances between the neighboring neurons. The lighter the colors between neurons, the closer the neurons are. The right

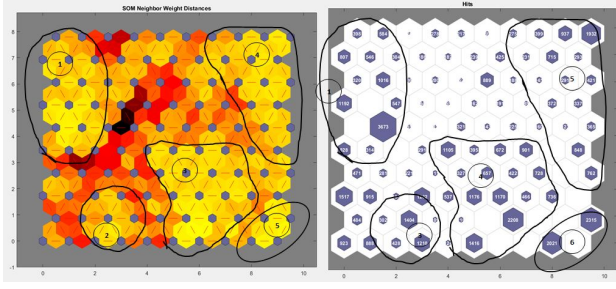


Fig. 3. Cluster Distribution of GTD data

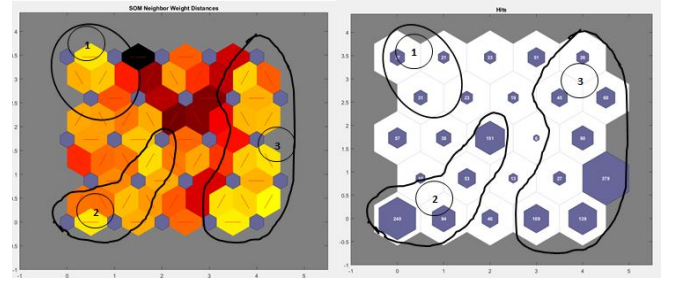


Fig. 5. Cluster Distribution of Twitter data with 'sentiment' feature

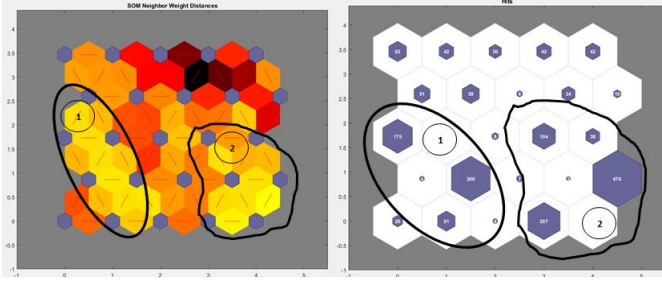


Fig. 4. Cluster Distribution of Twitter data

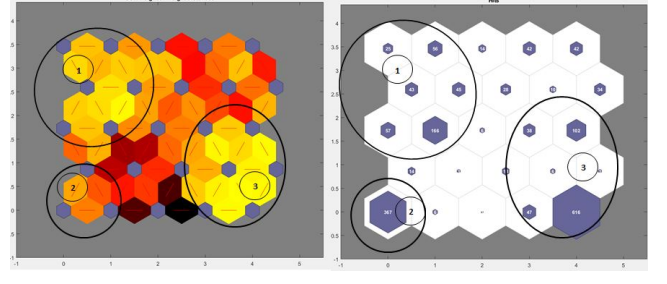


Fig. 6. Cluster Distribution of Twitter data with 'favorite' feature

side of figures are the hit histograms which represent the number of data instances that 'hit' to each neuron. The 'hit' neuron is the BMU of the input data instance.

The identified differences of the trends reflected by these four different data sets are summarized in Table II. The significant values of the features of each cluster were extracted based on the data that 'hit' to the neurons in the clusters.

Based on the clustering results and the values of the features for each cluster of Twitter data with four features, a cluster with moderately lower GTI value is identified. This cluster is not observed in clustering results of the GTD data. The region feature is one of the most significant features that characterizes the clusters learned by SOM using the Twitter data. The clusters learned by SOM on the GTD data focus on a few specific regions combined with values of GTI. To investigate whether Twitter data has any bias in countries with GTI values, we found that the average GTI of the user location was 5.24, and the average GTI of the country specific in the tweet was 5.51. There are 23% of tweets mentioning a specific country. The Twitter data does not show a bias in the regions with low GTI values.

Feature crit1, or 'motive for attack' is a significant feature that contributed to most of the clusters identified in the GTD data. However, it did not reflect in any of the clusters identified in the Twitter data sets. This indicates that motive for the attack is not a central similarity in the way Twitter users talk about terrorism.

In our database, 13% of the Twitter user mentioned a specific target type, which is larger than the number of hits received by most of the neurons in any cluster. The only target type that significantly contributed to the clustering results of Twitter data with four features is government type.

Those tweets 'hit' to neuron 15 in cluster 2, as shown in Figure 3. After further investigation into the Twitter data, we found that government type was the most commonly mentioned target type (70 tweets, 30% of the tweets with a specific target type). The second most common target type mentioned was private citizens/civilians type (51 tweets, 22% of the tweets with a specific target type) and the third most common instance was religious figures or institutions type (15 tweets, 6% of the tweets with a specific target type). This differed from the diverse clustering based on target type in the GTD; additionally, the GTD data primarily clustered around the target type of civilians and religious figures/institutions, demonstrating an over-representation of government focused attacks and a correlation with region 1 (North America) and mentioning government as a target type in the Twitter data.

The addition of two extra features, 'favorites' and 'sentiment', changed identified clusters on the Twitter data. This means these two features have a high impact on the clustering patterns. Adding the additional feature of user sentiment created three separate clusters from the Twitter data, as shown in 5. Particularly, an interesting trend observed is that a higher user sentiment is associated with region 6 (South Asia), and many of these tweets have the government as a target type. Adding 'favorite' as a feature also generated three clusters. It classified a United States and non-United States cluster with nonzero and zero favorites, respectively, suggesting that terrorism is discussed and sentiments are more readily agreed upon outside of the United States.

### C. Reflections of the Twitter Data

Given that countries, regions, and demographics use Twitter and social media disproportionately, not all public perception



TABLE II  
SUMMARY OF CLUSTERS AND VALUE OF THE FEATURES THAT ARE  
SIGNIFICANT TO THE CLUSTERS OF THE GTD AND TWITTER DATA

GTD data with four features	
Cluster No.	Significant Features and Values
1	Region: Middle East, Africa, Europe Target Type: private citizens religious figures/institutions, NGO
2	GTI: India, Afghanistan, Pakistan; Target Type: private citizens religious figures and institutions, terrorist groups, NGO
3	GTI: India, Afghanistan, Pakistan, Thailand Target Type: business, government, military, police
4	GTI: Middle East, Africa Target Type: Business, government, military, police
5	GTI: Nigeria/Syria Target Type: Police/military
Twitter Data with Four Features	
Cluster No.	Significant Features and Values
1	GTI: South Asia, Europe
2	GTI: US/Canada, Target Type: Government
Twitter data with 'sentiment' feature	
Cluster No.	Significant Features and Values
1	Low GTI: (average 1.67)
2	GTI: India, Pakistan, UK; Sentiment: Negative (Average -0.38)
3	GTI: US/Canada Sentiment: Neutral (Average -0.02)
Twitter data with 'favorite' feature	
Cluster No.	Significant Features and Values
1	GTI: Western Europe, Middle East, Africa
2	GTI: US/Canada Favorites: nonzero
3	GTI: Afghanistan, Pakistan, India (average: 7.66)

and sentiment is contained within Twitter. Since only one week of Twitter data is used in this research, the results from Twitter might be biased towards the most recent events. For example, during the one week tweets we used in this research, we identified two particular event are highly noted with the tweets: Canadian reparations to controversial figure Omar Khadr, and the US declaration of Pakistan as a terrorist state (which was widely discussed by Twitter users from India). These two events may have led to higher representation of Canada, India, and Pakistan in the Tweets surveyed.

## V. CONCLUSION AND FUTURE WORK

In this research work, we proposed to use unsupervised learning algorithm Self-Organizing Map to analyze and visualize the differences in trends and clusters identified in the Global Terrorist Database (GTD) data and Twitter data which was extracted with query words "terrorism" and "terrorist". We have also investigated the trends in Twitter data by adding the number of favorites and sentiment of the tweets as additional features. The clustering results showed that social media data clusters differently and weights input features differently than the objective GTD data. Additional features such as user favorites and user sentiment affected the way the data was clustered. The discrepancy in clustering between GTD and Twitter data may be explained by bias, disproportionate

representation, or change in terrorist trends not reflected in previous GTD records.

Future work includes extending this research to include data of a larger time span for a more comprehensive comparison and including additional user features, such as age, gender, demographics, and users followed to characterize people and their views on current events such as terrorism.

## ACKNOWLEDGMENT

This research was made possible with the support of the Indiana University-Purdue University Indianapolis Department of Computer and Information Technology, as well as through funding from the National Science Foundation and the United States Department of Defense. The authors would also like to thank Dr. Feng Li, Dr. Eugenia Fernandez, Sheila Walter, and Jovita Weah for their support.

## REFERENCES

- [1] J. Wagner and P. Rucker, Here are the 78 terrorist attacks the White House says were largely underreported, The Washington Post, 06-Feb-2017. [Online]. Available: <https://www.washingtonpost.com/news/post-politics/wp/2017/02/06/here-are-the-78-terrorist-attacks-the-white-house-says-were-largely-under-reported/>. [Accessed: 28-Jul-2017].
- [2] K. Zaman, ISIS Has a Twitter Strategy and It Is Terrifying [Infographic], Medium, 20-Nov-2015. [Online]. Available: <https://medium.com/fifth-tribe-stories/isis-has-a-twitter-strategy-and-it-is-terrifying-7cc059ccf51b.fg5420c2v>. [Accessed: 28-Jul-2017].
- [3] J. Constine, Facebook requests input on hard questions about censorship and terrorism, TechCrunch, 15-Jun-2017. [Online]. Available: <https://techcrunch.com/2017/06/15/facebook-censorship-terrorism/>. [Accessed: 28-Jul-2017].
- [4] Xi Meng, Haowen Mo, Shenhe Zhao and Jianqiang Li, "Application of anomaly detection for detecting anomalous records of terrorist attacks," 2017 IEEE 2nd International Conference on Cloud Computing and Big Data Analysis (ICCCBDA), Chengdu, 2017, pp. 70-75.
- [5] P. Misra, Development and Optimization of Machine Learning Algorithms and Models of Relevance to START Databases, National Consortium for the Study of Terrorism and Responses to Terrorism, rep., Apr. 2016.
- [6] N. Azam, Jahiruddin, M. Abulaish and N. A. H. Haldar, "Twitter Data Mining for Events Classification and Analysis," 2015 Second International Conference on Soft Computing and Machine Intelligence (ISCMI), Hong Kong, 2015, pp. 79-83.
- [7] K. Byrd, A. Mansurov and O. Baysal, "Mining Twitter Data for Influenza Detection and Surveillance," 2016 IEEE/ACM International Workshop on Software Engineering in Healthcare Systems (SEHS), Austin, TX, 2016, pp. 43-49. doi: 10.1109/SEHS.2016.016
- [8] L. Montesinos, S. J. P. Rodriguez, M. Orchard and S. Eyheramendy, "Sentiment analysis and prediction of events in TWITTER," 2015 CHILEAN Conference on Electrical, Electronics Engineering, Information and Communication Technologies (CHILECON), Santiago, 2015, pp. 903-910.
- [9] O. Oh, M. Agrawal, and H. R. Rao, Information control and terrorism: Tracking the Mumbai terrorist attack through twitter, Information System Frontiers, vol. 13, no. 1, pp. 3343, Mar. 2011.
- [10] M. Cheong and V. C. S. Lee, A microblogging-based approach to terrorism informatics: Exploration and chronicling civilian sentiment and response to terrorism events via Twitter, Information Systems Frontiers, vol. 13, no. 1, pp. 4549, Mar. 2011.
- [11] T. Kohonen, Self-organizing maps. Berlin: Springer, 1995.
- [12] National Consortium for the Study of Terrorism and Responses to Terrorism (START). (2017). Global Terrorism Database [Data file]. Retrieved from <https://www.start.umd.edu/gtd>
- [13] economicsandpeace.org, economicsandpeace.org. Institute for Economics and Peace, Nov-2016.
- [14] REST API, Twitter Developer Documentation, 2017. [Online]. Available: <https://dev.twitter.com/rest/public>. [Accessed: 31-Jul-2017].
- [15] R. Petersens, AFINN, 2017. [Online]. Available: <http://www2.imm.dtu.dk/pubdb/views/publicationdetails.php?id=6010>. Accessed: 31-Jul-2017